

ERRORS IN VARIABLES AND SERIALLY CORRELATED
DISTURBANCES IN DISTRIBUTED LAG MODELS¹

BY D. M. GREYER AND G. S. MADDALA

I. INTRODUCTION

FOR SOME TIME econometricians have been bothered by the very long lags in response that their estimates often show. There have been some attempts to explain these apparent long lags as statistical artifacts due to serial correlation bias, time aggregation, misspecification of the lag distribution, etc. It is well known that errors in variables cause least squares estimates of regression coefficients to be biased and inconsistent, but there does not seem to be any full treatment in the econometric literature of the problem of errors in variables in the context of distributed lag models.

In this paper we consider three types of models:

$$(i) \quad y_t = \gamma_0 x_t + \gamma_1 x_{t-1} + \dots + \gamma_k x_{t-k} + u_t,$$

$$(ii) \quad y_t = \alpha \sum_{i=0}^{\infty} \delta^i x_{t-i} + v_t,$$

$$(iii) \quad y_t = \alpha y_{t-1} + \gamma x_t + w_t.$$

In model (i) y_t depends only upon the current level of x_t and possibly a finite number of lagged x 's as well. Model (ii) is a distributed lag model in distributed lag form; i.e., the model is specified with only exogenous variables on the right hand side. Model (iii) is a distributed lag model in autoregressive form; i.e., lagged values of the dependent variable are included among the set of explanatory variables. In each case we assume that the exogenous variables are observed with a measurement error η_t which is uncorrelated with all other variables in the model. To the extent that the measurement errors are sampling errors it seems reasonable to expect them to be uncorrelated over time, though in most of what follows we do not impose this restriction. Thus, we do allow for the possibility of some systematic observation errors possibly due to infrequent revision of sampling procedures or to persistent errors in estimating subjective variables such as expectations. Throughout this paper we assume that estimates of the structural parameters are desired. Thus, we take the relevant question to be: given a change in x_t , what is the magnitude and timing of the response in y ? This question would be important to a policy maker who could influence the future levels of x_t . We do not consider the problem of how to predict y_t given observations on it and the

¹ This research was supported by grants from the National Science Foundation and from the Ford Foundation to the Cowles Foundation for Research in Economics at Yale. We thank Professors Goldberger and Kmenta and the referees for helpful comments. Responsibility for any errors is solely ours.

imperfectly measured x 's. In Section 2 we examine models with no lagged endogenous variables. Section 3 treats distributed lag models in the autoregressive form. The conclusions of the paper are given in Section 4.

In what follows all results will be given in terms of the population moments of the variables involved. Strictly speaking the equations should be expressed as probability limits, but for economy of notation the plim will be omitted. In all cases it is assumed that the measurement errors are uncorrelated with other variables in the model and that all variables have mean zero. The following notational conventions will be observed throughout:

$$\begin{aligned}\text{cov}(x_t, y_{t-j}) &= \sigma_{xy}(j), & j \neq 0, \\ &= \sigma_{xy}, & j = 0, \\ \frac{\sigma_{xy}(j)}{\sigma_x \sigma_y} &= \rho_{xy}(j).\end{aligned}$$

2. ERRORS IN VARIABLES IN MODELS WITHOUT LAGGED ENDOGENOUS VARIABLES

To begin with, suppose one has data on y_t where:

$$\begin{aligned}(1.1) \quad y_t &= \gamma x_t + u_t, & E(u_t) &= E(x_t) = 0, & E(u_t x_{t-j}) &= 0, & \forall j; \\ z_t &= x_t - \eta_t, & E(\eta_t) &= 0, & E(\eta_t x_{t-j}) &= 0, & \forall j, \\ & & & & E(\eta_t u_{t-j}) &= 0, & \forall j.\end{aligned}$$

In this case it is well known that the least squares estimate of γ obtained by regressing y on z has the limit

$$(1.2) \quad \hat{\gamma} = \gamma(1 - \lambda) \quad \text{where} \quad \lambda = \frac{\sigma_\eta^2}{\sigma_x^2 + \sigma_\eta^2}.$$

Thus, one has the usual result that at least for models with one independent variable, measurement errors lead to estimates that are too close to zero.

If e_t is the calculated residual from the regression, then

$$\begin{aligned}(1.3) \quad e_t &= y_t - \hat{y}_t = \gamma x_t + u_t - \hat{\gamma} z_t \\ &= u_t + \gamma \lambda x_t - \gamma(1 - \lambda)\eta_t.\end{aligned}$$

The first order autocorrelation of e_t is equal to

$$(1.4) \quad \rho_{ee}(1) = \frac{\gamma^2 \lambda^2 \sigma_{xx}(1) + \gamma^2 (1 - \lambda)^2 \sigma_{\eta\eta}(1) + \sigma_{uu}(1)}{\gamma^2 \lambda^2 \sigma_x^2 + \gamma^2 (1 - \lambda)^2 \sigma_\eta^2 + \sigma_u^2}.$$

Thus, the estimated serial correlation is a weighted average of $\rho_{xx}(1)$, $\rho_{\eta\eta}(1)$, and $\rho_{uu}(1)$ with weights $\gamma^2 \lambda^2 \sigma_x^2$, $\gamma^2 (1 - \lambda)^2 \sigma_\eta^2$, and σ_u^2 , respectively. Note that if there is no serial correlation in u , then the presence of measurement errors will produce

calculated residuals which are autocorrelated.² Though, if the measurement errors are independent (a not unreasonable assumption), then the autocorrelation in the residuals is likely to be quite small. For instance, if ρ_{xy}^2 is 0.8 and 10 per cent of the variance of z is due to η , then $\rho_{ee}(1)$ is only $\rho_{xx}(1)/35$. Even if half the variance of z comes from measurement errors, $\rho_{ee}(1)$ is only $\rho_{xx}(1)/3$. To the extent that autocorrelated residuals may be taken as signifying missing variables or some kind of distributed lag model there seems little chance that serially independent errors of measurement would mislead an investigator into mistakenly specifying a model with lagged values of y or z . The danger is much greater, however, if the measurement errors are correlated across observations.

When serial correlation is found in the calculated residuals it is common practice to reestimate the model using the ρ -differenced data. If γ is estimated by regressing $y_t - \rho y_{t-1}$ on $z_t - \rho z_{t-1}$ where ρ is some number between plus and minus one, the resulting estimate has the limit

$$(1.5) \quad \hat{\gamma} = \frac{\gamma \sigma_x^2 (1 + \rho^2 - 2\rho \rho_{xx}(1))}{\sigma_x^2 (1 + \rho^2 - 2\rho \rho_{xx}(1)) + \sigma_\eta^2 (1 + \rho^2 - 2\rho \rho_{\eta\eta}(1))}.$$

Comparing (1.5) with (1.2), note that if ρ is positive $|\hat{\gamma} - \gamma|$ is increased whenever $\rho_{xx}(1)$ is greater than $\rho_{\eta\eta}(1)$, and vice versa for ρ negative. Since positive autocorrelation is generally the case with economic time series and since it is reasonable to expect that x_t is more strongly serially correlated than η_t , transformations often employed to increase efficiency may well increase $|\hat{\gamma} - \gamma|$ (see, also, Sims [4]).

Suppose that one suspects that past levels of x as well as current x enter into the determination of y and estimates the following model:

$$(1.6) \quad y_t = g z_t + h z_{t-n} + w_t,$$

whereas the true model is still (1.1). The least squares coefficients have the limits

$$\begin{aligned}(1.7) \quad \hat{g} &= \frac{\gamma(\sigma_x^4(1 - \rho_{xx}^2(n)) + \sigma_\eta^2 \sigma_x^2(1 - \rho_{xx}(n)\rho_{\eta\eta}(n)))}{D}, \\ \hat{h} &= \frac{\gamma \sigma_x^2 \sigma_\eta^2 (\rho_{xx}(n) - \rho_{\eta\eta}(n))}{D},\end{aligned}$$

where $D = \sigma_x^4(1 - \rho_{xx}^2(n))$. Equations (1.7) show that the presence of measurement errors will result in non-zero coefficients of lagged z 's even though the true model contains no lags at all. Note that this is not a consequence of serially correlated measurement errors, but arises whenever the autocorrelation functions of x and η_t differ. If y depends upon $\{x_{t-j}, j = 0, 1, \dots, k\}$, then it will normally be the case that z 's lagged more than k periods will have non-zero coefficients in the regression

² Except in the case where both x and η are serially independent. Since most economic time series are highly autocorrelated it is not likely that $\rho_{xx}(1)$ is zero. If $\rho_{xx}(1)$ and $\rho_{\eta\eta}(1)$ are zero, then the calculated residuals will be less strongly autocorrelated than the sequence $\{u_t\}$.

of y_t on $\{z_{t-j}, j = 0, 1, \dots, k, k+1, \dots, m\}$. Thus, measurement errors may result in the appearance of adjustment lags which are longer than the correct lag.³

Though adding additional lagged z 's to the regression gives a misleading picture of the shape of the lag distribution, it may provide a better estimate of the total response. While the sum of the coefficients in (1.7) is necessarily closer to zero than γ , the sum will be a better estimate of γ than that given in equation (1.2) provided that x is more strongly correlated than η . Since this condition is likely to be satisfied in practice, the obvious conclusion is that adding lagged z to the regression will reduce the inconsistency in γ , but one should not infer anything else from the implied lag.

Consider now the case in which the true model is a distributed lag model in distributed lag form. Suppose that the distributed lag is of the Koyck type:

$$(1.8) \quad y_t = \frac{\alpha}{1 - \delta L} x_t + e_t$$

where L is the lag operator defined by $Lx_t \equiv x_{t-1}$. Due to measurement errors we estimate

$$(1.9) \quad y_t = \frac{\alpha z_t}{1 - \delta L} + w_t$$

where $z_t = x_t + \eta_t$. Define

$$(1.10) \quad X_t^*(\delta) = \frac{x_t}{1 - \delta L}, \quad \eta_t^*(\delta) = \frac{\eta_t}{1 - \delta L}, \quad Z_t^*(\delta) = X_t^*(\delta) + \eta_t^*(\delta),$$

where $|\delta| < 1$.

The maximum likelihood procedure for estimating α and δ (with no measurement errors) is as follows:⁴ For values of δ less than one in absolute value, compute X_t^* and regress y_t on X_t^* to obtain an estimate of α . Choose as estimates the values $(\hat{\alpha}, \hat{\delta})$ which minimize the residual variance. If, due to the presence of measurement errors, we use Z_t^* instead of X_t^* , what are the effects on the estimates of α and δ ?

First, consider the case in which both x_t and η_t are uncorrelated over time. In this case

$$(1.11) \quad \begin{aligned} \text{var}(X_t^*(\delta)) &= \frac{\sigma_x^2}{1 - \delta^2}, \\ \text{var}(\eta_t^*(\delta)) &= \frac{\sigma_\eta^2}{1 - \delta^2}, \\ \hat{\alpha}(\delta) &= \alpha \left(\frac{1 - \delta^2}{1 - \delta\bar{\delta}} \right) (1 - \lambda). \end{aligned}$$

³ This kind of result easily generalizes by the type of argument Theil [5] applied to the analysis of specification errors. In general any z 's that have non-zero coefficients in the projection of x_t on $\{z_t\}$ will have non-zero estimated coefficients in a regression explaining y_t . Note that if y_t depends only on z_{t-n} , then z 's lagged less than n periods will also appear in the regression, implying too short a lag.

⁴ See Klein [2].

Since the ratio of the variance of X^* to the variance of η^* is independent of δ , the estimate of α is too close to zero by a constant proportion for each choice of δ . Similarly, the explained variance is too small by a constant factor for each δ . Let \hat{u}_t be the calculated residuals; then

$$(1.12) \quad \text{var}(\hat{u}) = \sigma_y^2 - \alpha^2 \sigma_x^2 \frac{(1 - \delta^2)}{(1 - \delta\bar{\delta})^2} (1 - \lambda).$$

So while $|\hat{\alpha}|$ is less than $|\alpha|$, the value of δ which minimizes the residual variance is unaffected by the measurement errors. The situation changes, however, if it is assumed that x_t is serially correlated. If x_t is a first order autoregressive process with parameter ρ , then

$$(1.13) \quad \begin{aligned} \text{var}(X_t^*(\delta)) &= \frac{\sigma_x^2(1 + \rho\delta)}{(1 - \bar{\delta}\rho)(1 - \bar{\delta}^2)}, \\ \text{var}(\eta_t^*(\delta)) &= \frac{\sigma_\eta^2}{(1 - \bar{\delta}^2)}, \\ \hat{\alpha}(\delta) &= \alpha \frac{(1 - \delta^2)(1 - \rho^2\delta\bar{\delta})}{(1 - \delta\bar{\delta})(1 - \rho\delta)(1 + \rho\bar{\delta})} \left\{ \frac{\sigma_x^2(1 + \rho\delta)}{\sigma_x^2(1 + \rho\delta) + \sigma_\eta^2(1 - \rho\bar{\delta})} \right\}, \end{aligned}$$

and

$$\text{var}(\hat{u}) = \sigma_y^2 - \frac{\alpha^2(1 - \delta^2)(1 - \rho^2\delta\bar{\delta})\sigma_x^2}{(1 - \delta\bar{\delta})^2(1 - \rho\delta)(1 - \rho^2\bar{\delta}^2)} \left\{ \frac{\sigma_x^2(1 + \rho\delta)}{\sigma_x^2(1 + \rho\delta) + \sigma_\eta^2(1 - \rho\bar{\delta})} \right\},$$

where the terms in curly brackets show the effects of the measurement errors. The effects of errors in variables are no longer independent of δ . The derivative of $\sigma_x^2(1 + \rho\delta)/(\sigma_x^2(1 + \rho\delta) + \sigma_\eta^2(1 - \rho\bar{\delta}))$ with respect to δ has the same sign as ρ , so if ρ is positive, as is likely to be the case, the estimate of $|\alpha|$ is too small, and the estimate of δ is too large, giving the result that errors of measurement will make the adjustment appear slower than is actually the case. These results hold for any stationary specification of the disturbances u_t .

Continuing with this example, we can ask about the nature of the apparent distributed lag relationship between y and z . Consider the projection of y_t on the current past values of z_t , that is, the lag distribution γ such that

$$(1.14) \quad y_t = \gamma(L)z_t + \phi_t$$

where $E(\phi_t z_{t-j}) = 0$, $j \geq 0$. A straightforward but somewhat tedious application of the results given by Whittle [6] gives that $\gamma(L)$ is of the form $(A + BL)/(1 - \beta L)$ $(1 - \delta L)$.⁵ Thus the measurement errors lead to the appearance of a higher order distributed lag than is actually the case.

⁵ See Whittle [6], especially Example 3.3.7, pp. 35, and Ch. 8, Theorem 1, p. 93. The expression for β is shown on p. 35. Assuming that δ and ρ are positive, then the mean lag implied in (1.14) is greater than the true mean lag. Knowledge of the parameters in (1.14) does not allow one to determine the structural parameters without knowing the true models for y_t , x_t , and η_t .

3. MODELS WITH LAGGED ENDOGENOUS VARIABLES AND SERIALLY CORRELATED ERRORS

In this section we discuss the effects of measurement errors when the estimating equation is of the form:

$$(2.1) \quad \hat{y}_t = \hat{\alpha}y_{t-1} + \hat{\gamma}z_t.$$

Consider, initially, the case in which the true model is that given in (1.1), i.e., the true coefficient of y_{t-1} is zero. In this case the least squares estimates have the limits

$$(2.2) \quad \begin{aligned} \hat{\alpha} &= \frac{\sigma_\eta^2 \sigma_y^2 \rho_{yy}(1)}{D} + \frac{\sigma_x^2 \sigma_{uu}(1)}{D}, \\ \hat{\gamma} &= \frac{\gamma(D - \sigma_\eta^2 \sigma_y^2)}{D} - \frac{\gamma \sigma_{xx}(1) \sigma_{uu}(1)}{D}, \end{aligned}$$

where $D = \sigma_y^2 \sigma_z^2 - \sigma_{xy}^2(1)$. If the x sequence is positively correlated over time and the disturbances are either independent or positively correlated, then $\hat{\alpha}$ will be positive, and the presence of measurement errors will cause $\hat{\alpha}$ to be larger than would otherwise be the case.⁶ In any case $|\hat{\alpha}|$ is less than one so that the estimate obtained will most likely appear reasonable.

If the true coefficient of y_{t-1} is not zero, then the least squares estimates have the limits

$$(2.3) \quad \begin{aligned} \hat{\alpha} - \alpha &= \frac{\sigma_\eta^2 \sigma_y^2 (\rho_{yy}(1) - \alpha)}{D} + \frac{\sigma_x^2 \sigma_{uy}(1)}{D}, \\ \hat{\gamma} - \gamma &= \frac{-\gamma \sigma_\eta^2 \sigma_y^2}{D} - \frac{\gamma \sigma_{xy}(1) \sigma_{uy}(1)}{D}, \end{aligned}$$

$$D = \sigma_y^2 \sigma_z^2 - \sigma_{xy}^2(1).$$

Adding the assumptions that x_t and u_t are first order autoregressions with parameters ρ_x and ρ_u respectively, these expressions become:

$$(2.4) \quad \begin{aligned} (\hat{\alpha} - \alpha) &= -\frac{\sigma_\eta^2 \left(\frac{\rho_x \gamma^2 \sigma_x^2}{1 - \alpha \rho_x} + \frac{\rho_u \sigma_u^2}{1 - \alpha \rho_u} \right)}{D} + \frac{\rho_u \sigma_u^2 \sigma_x^2}{D(1 - \alpha \rho_u)}, \\ (\hat{\gamma} - \gamma) &= \frac{-\gamma \sigma_\eta^2 \sigma_y^2}{D} - \frac{1}{D} \cdot \frac{\gamma \rho_x \sigma_x^2 \rho_u \sigma_u^2}{(1 - \alpha \rho_x)(1 - \alpha \rho_u)}. \end{aligned}$$

⁶ If the disturbances are negatively correlated, the bias in $\hat{\alpha}$ will be decreased.

These expressions show the inconsistency as a sum of two effects; the first term is due to measurement errors only, and the second (except for a slight modification of D) is the effect of serial correlation.⁷

Serial correlation will lead to an overestimate of α provided ρ_u is positive, and if ρ_x is also positive, errors of measurement will have a similar effect. The presence of errors in variables always leads to an estimate of γ which is too close to zero, and this effect will be increased or decreased as ρ_x and ρ_u are of the same or different signs. Thus, in the most frequently assumed case (ρ_x and ρ_u both positive), the two effects work in the same direction to make the inconsistency greater.

4. CONCLUSIONS

In this paper we have examined the effects of errors in variables on the probability limits of the estimated coefficients in distributed lag models. In the case of models with no lags or with finite lags, measurement errors in the exogenous variable may lead to the appearance of spurious long lags in adjustment. We have also shown that measurement errors may produce residuals which are autocorrelated and that methods frequently used to increase efficiency when the true disturbances are serially correlated are likely to result in increased inconsistency.

For distributed lag models estimated in the distributed lag form we have shown that under plausible conditions measurement errors may lead to estimates which imply that adjustment is slower than is actually the case. Finally, we have derived expressions for the effects of serial correlation and errors in variables for a model with a single exogenous variable and a lagged endogenous variable. In this case also it is likely that the presence of measurement errors will augment the effects of serial correlation and give estimates of rates of adjustment which are too slow.

*California Institute of Technology
and
University of Rochester*

Manuscript received May, 1971; revision received October, 1971.

⁷ These expressions generalize those given by Griliches [1] and Malinvaud [3]. Note that if x and η are known to be first order autoregressive processes, then from the second moments of z and the covariances between y_t and z_{t-j} it is possible to obtain consistent estimates of ρ_x , ρ_u , σ_x^2 , and σ_η^2 . Thus, provided one is willing to invest in a complete model for the variables involved, all the parameters may be consistently estimated.

To be explicit, consider the model in equations (i). If, in addition, we make the assumption that x_t and η_t are first order autoregressive, then it is easy to verify that

$$\text{cov}(z_t, z_{t-j}) = \rho_x^j \sigma_x^2 + \rho_\eta^j \sigma_\eta^2, \quad \rho_x \neq \rho_\eta \quad (j = 0, 1, 2, 3).$$

These four equations can be solved to get consistent estimates of ρ_x , ρ_η , σ_x^2 , and σ_η^2 , and thus, one can correct the least squares estimate of γ for the (asymptotic) bias. In the case of the distributed lag model (iii), one would have to exploit the covariances between y_t and lagged values of y_t also. The important point to note is that if everything in the bias term is estimable, then we can eliminate the bias, and we can do this if we add enough assumptions to the structure of the model. In actual practice, however, it is doubtful how much mileage one can get by this procedure.

REFERENCES

- [1] GRILICHES, ZVI: "Distributed Lags: A Survey," *Econometrica*, 34 (1967), 16-49.
- [2] KLEIN, L.: "The Estimation of Distributed Lags," *Econometrica*, 26 (1958), 553-565.
- [3] MALINVAUD, E.: *Statistical Methods of Econometrics*. Chicago: Rand McNally, 1966.
- [4] SIMS, C. A.: "The Role of Approximate Prior Restrictions in Distributed Lag Estimation," *Journal of the American Statistical Association*, 67 (1972), 164-175.
- [5] THEIL, H.: *Economic Forecasts and Policy*. Amsterdam: North-Holland, 1961.
- [6] WHITTLE, P.: *Prediction and Regulation by Linear Least Squares Methods*. London: The English Universities Press, 1963.